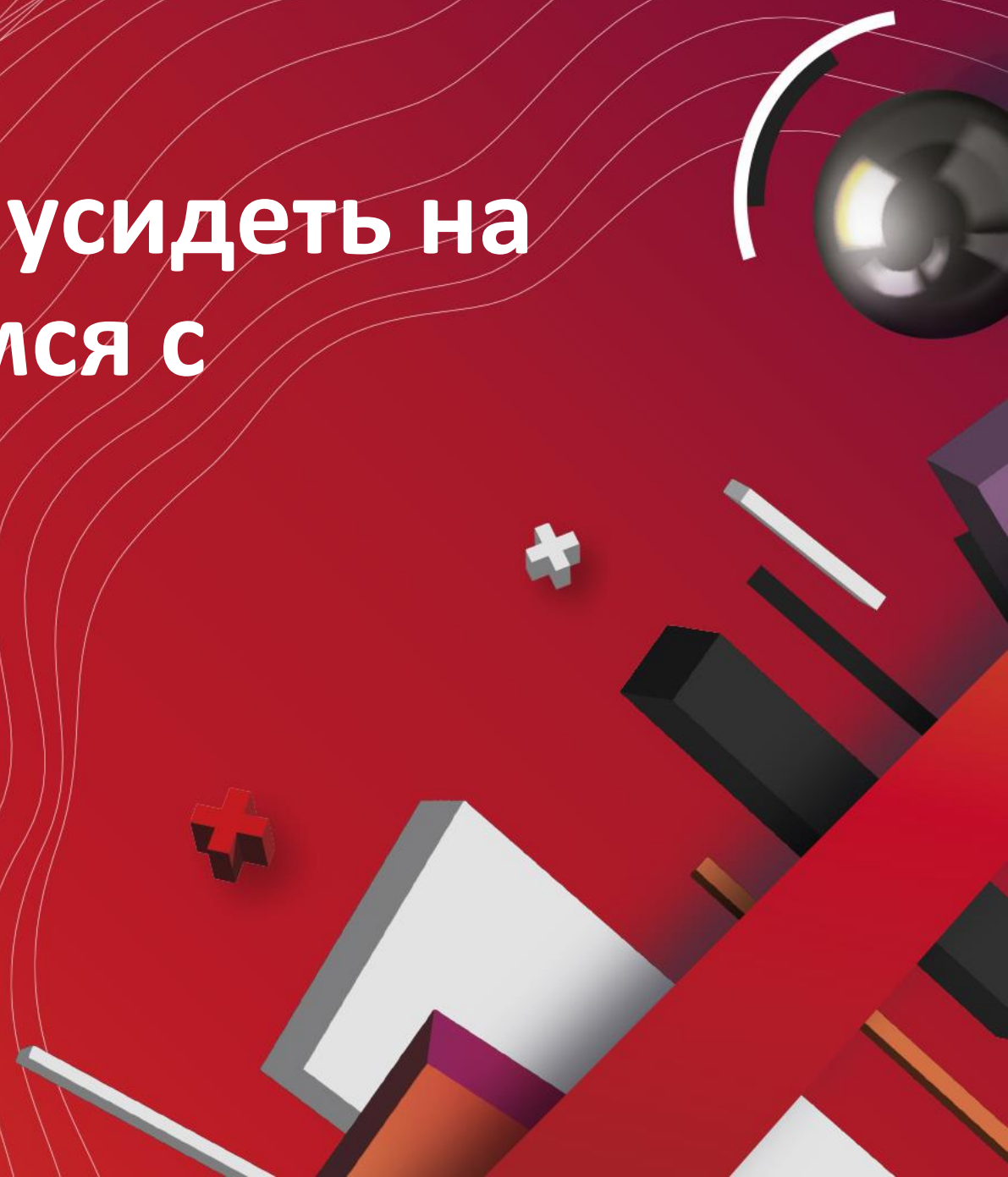


сАР и СаР: пытаемся усидеть на двух стульях и боремся с последствиями

Сергей Петренко



HighLoad⁺⁺
2022



О себе

Меня зовут Сергей Петренко

О себе

Меня зовут Сергей Петренко

Занимаюсь репликацией в Tarantool

О себе

Меня зовут Сергей Петренко

Занимаюсь репликацией в Tarantool

Занимался алгоритмом синхронной репликации
и выборов лидера

План

- Введение
- Особенности Raft в Tarantool
- Проблема асинхронных транзакций
- Размещение кластера Raft в двух ЦОДах
- Обнаружение различных историй лидерства
- Выводы

CAP теорема

В присутствии проблем сети (P) система может гарантировать либо консистентность (C), либо доступность (A), **но не то и другое вместе**

CAP теорема

В присутствии проблем сети (P) система может гарантировать либо консистентность (C), либо доступность (A), **но не то и другое вместе**

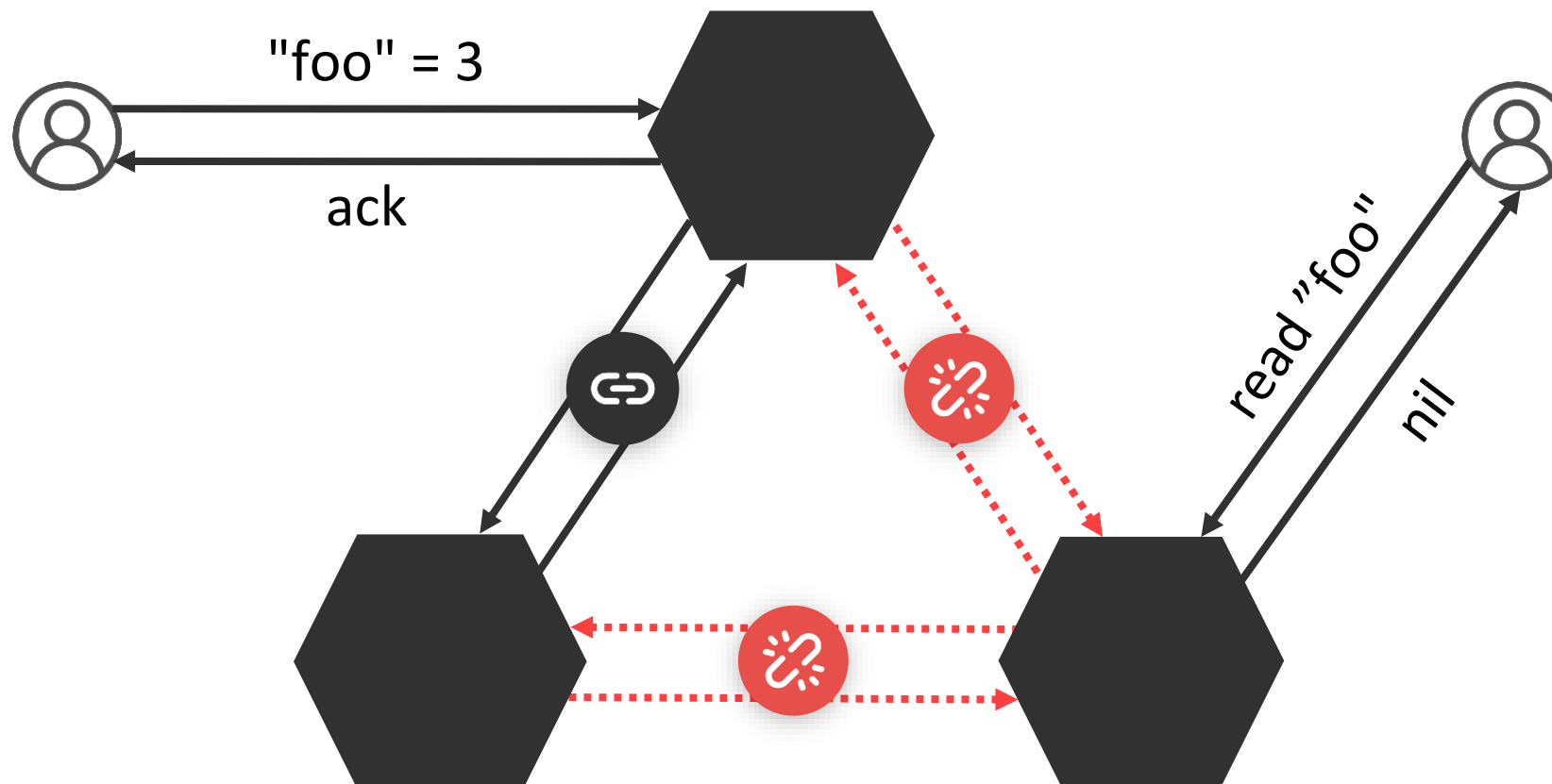
- Консистентность – запрос на чтение возвращает все записи, подтверждённые до этого запроса

CAP теорема

В присутствии проблем сети (P) система может гарантировать либо консистентность (C), либо доступность (A), **но не то и другое вместе**

- Консистентность – запрос на чтение возвращает все записи, подтверждённые до этого запроса
- Доступность – запрос к любому из работающих узлов возвращает ответ за конечное время

Почему бы не выбрать CAP?



Виды репликации

Асинхронная



Синхронная



Виды репликации

Асинхронная

- Available



Синхронная



Виды репликации

Асинхронная

- Available



Синхронная

- He available



Виды репликации

Асинхронная



- Available

Синхронная



- Не available
- Может быть consistent

Raft

- CaP
- кворум – большинство

План

- Введение
- Особенности Raft в Tarantool
- Проблема асинхронных транзакций
- Размещение кластера Raft в двух ЦОДах
- Обнаружение различных историй лидерства
- Выводы

В Raft нам не хватает режима cAP (вместо CaP)

Особенности нашей реализации Raft

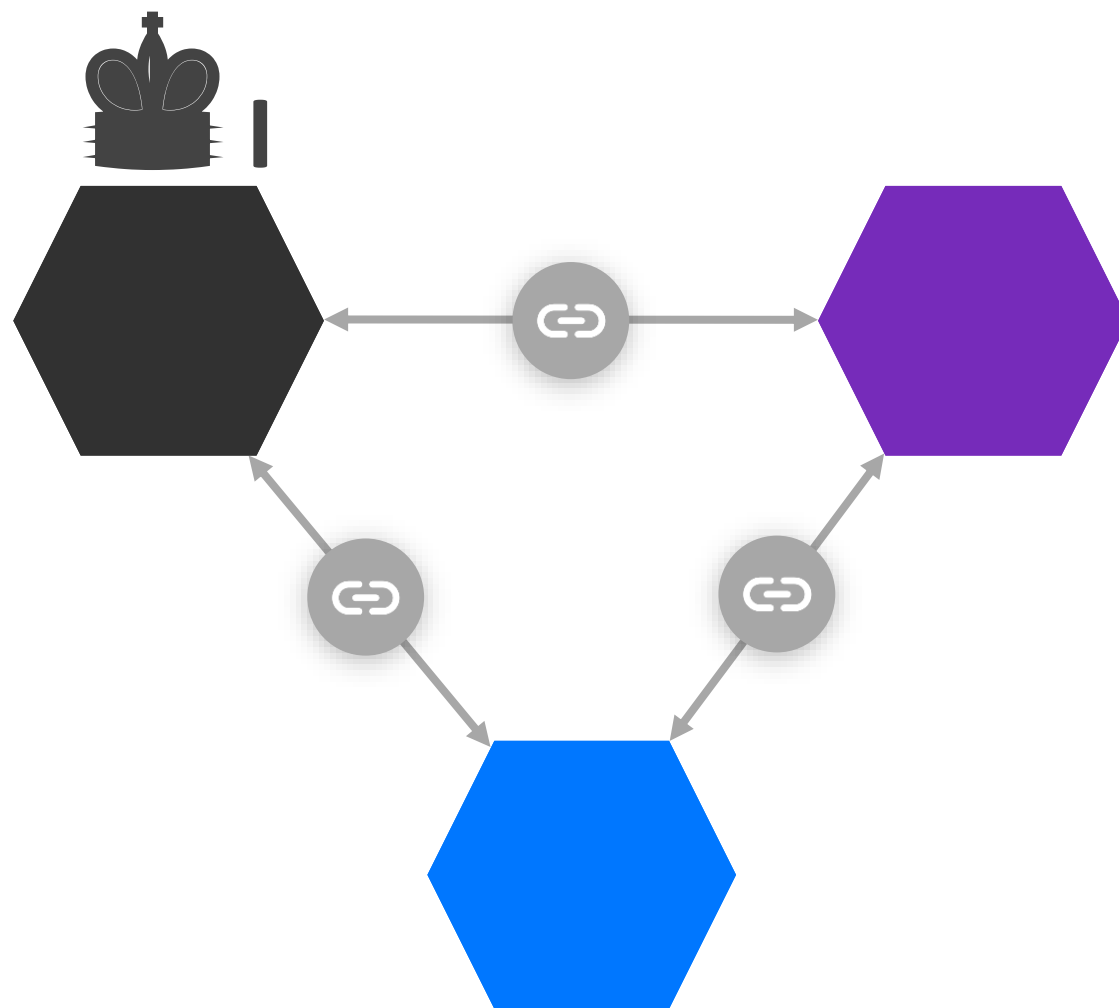
- Выбор между синхронной и асинхронной репликацией

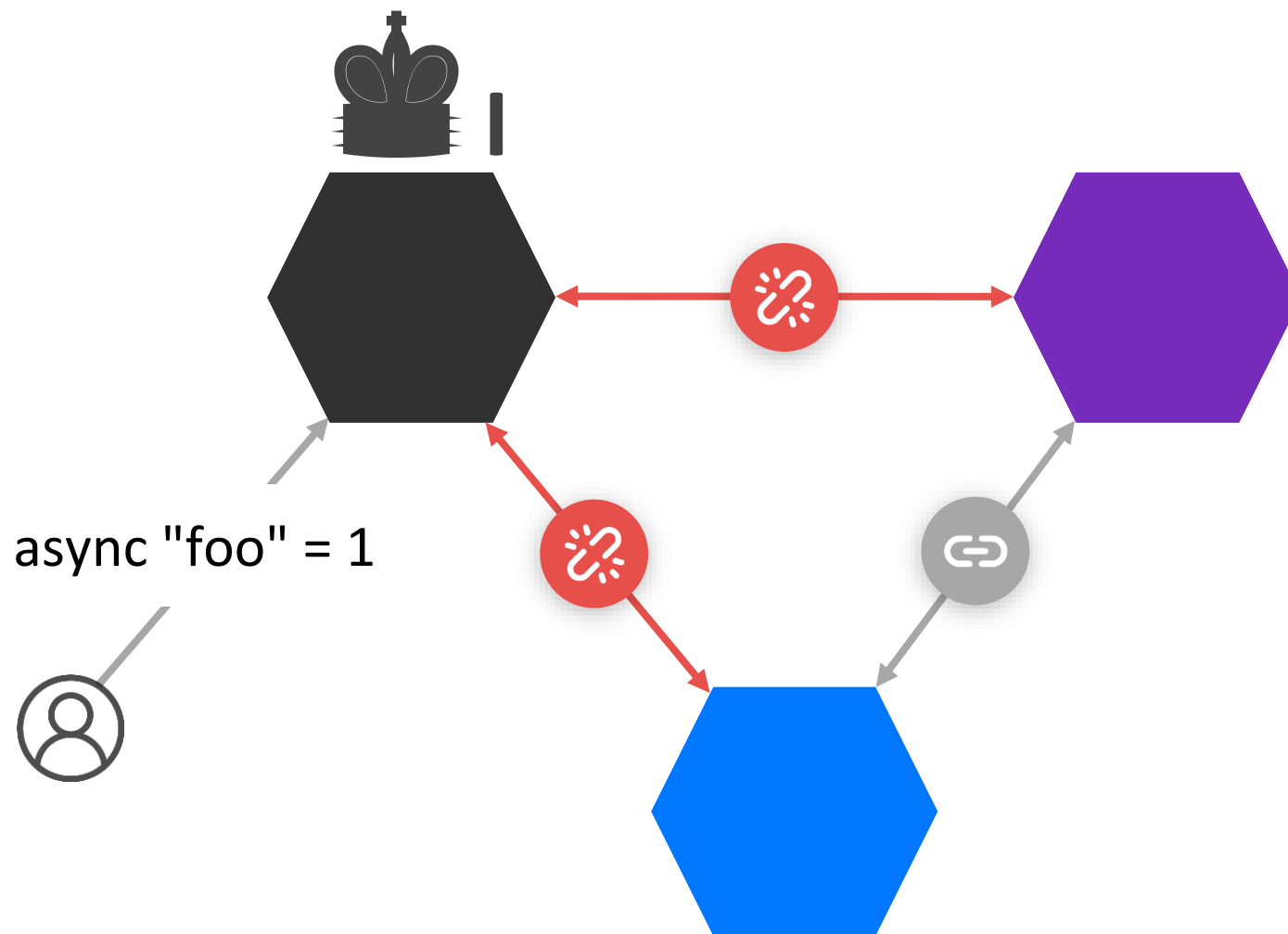
Особенности нашей реализации Raft

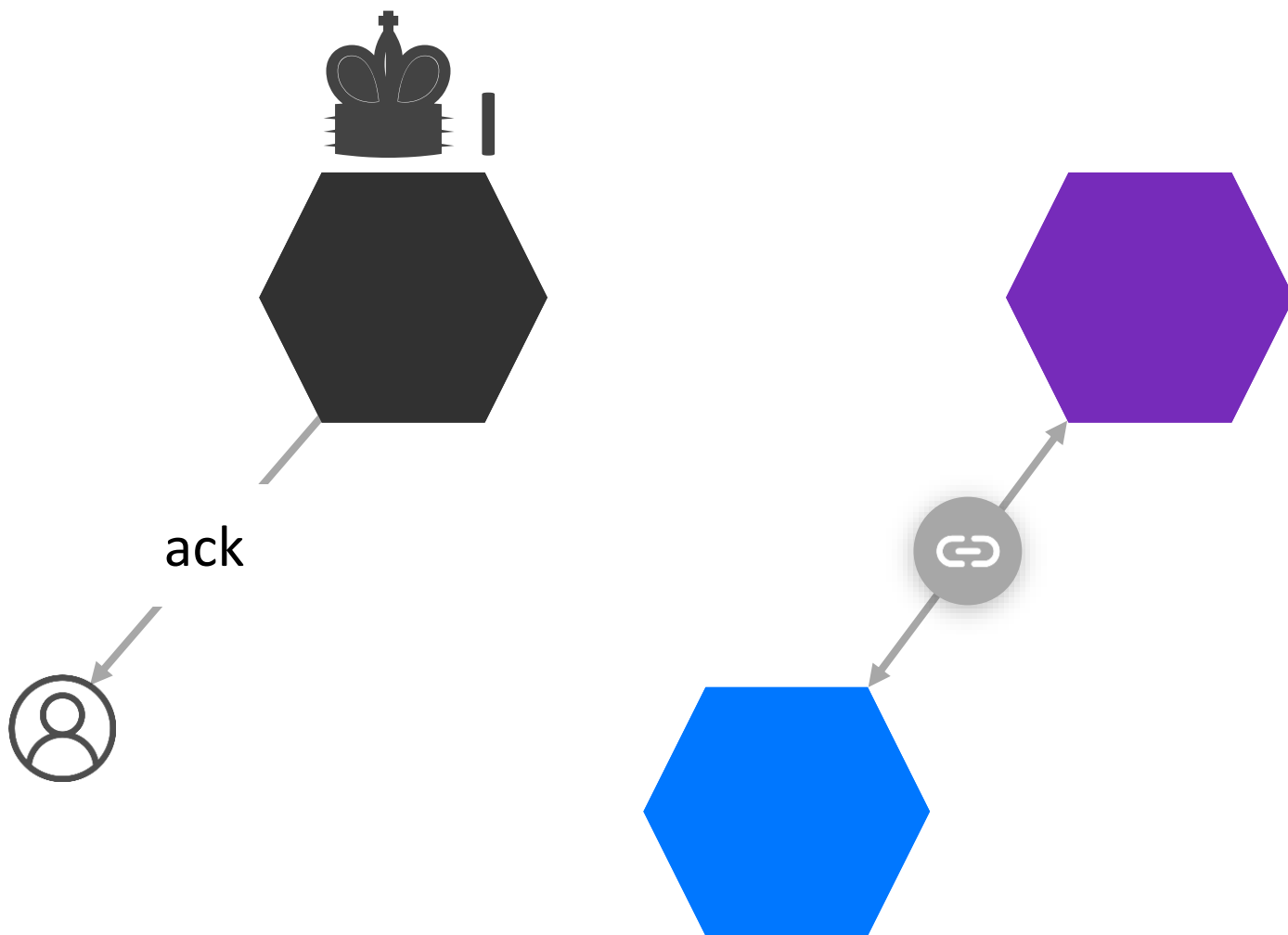
- Выбор между синхронной и асинхронной репликацией
- Настраиваемый кворум

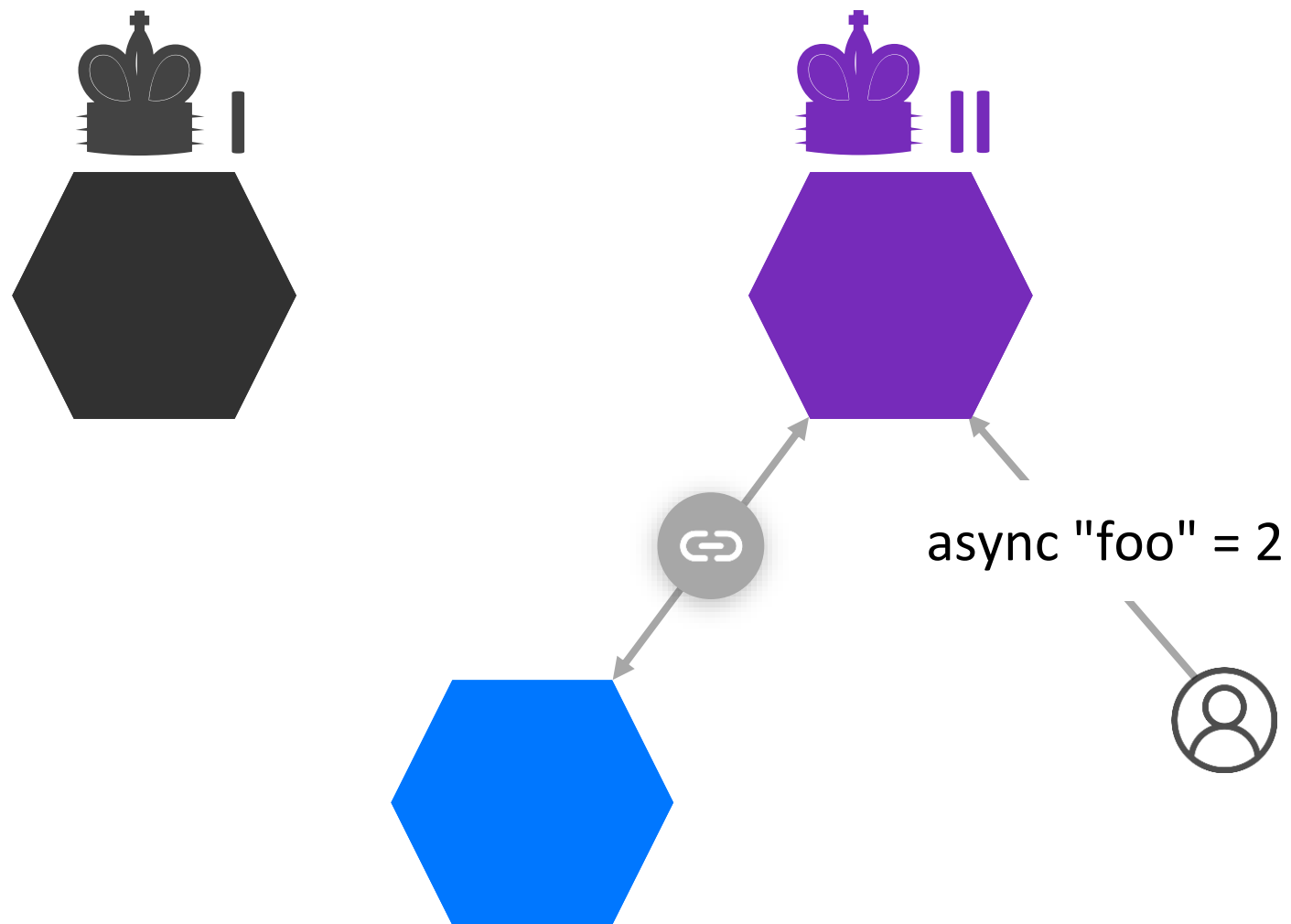
План

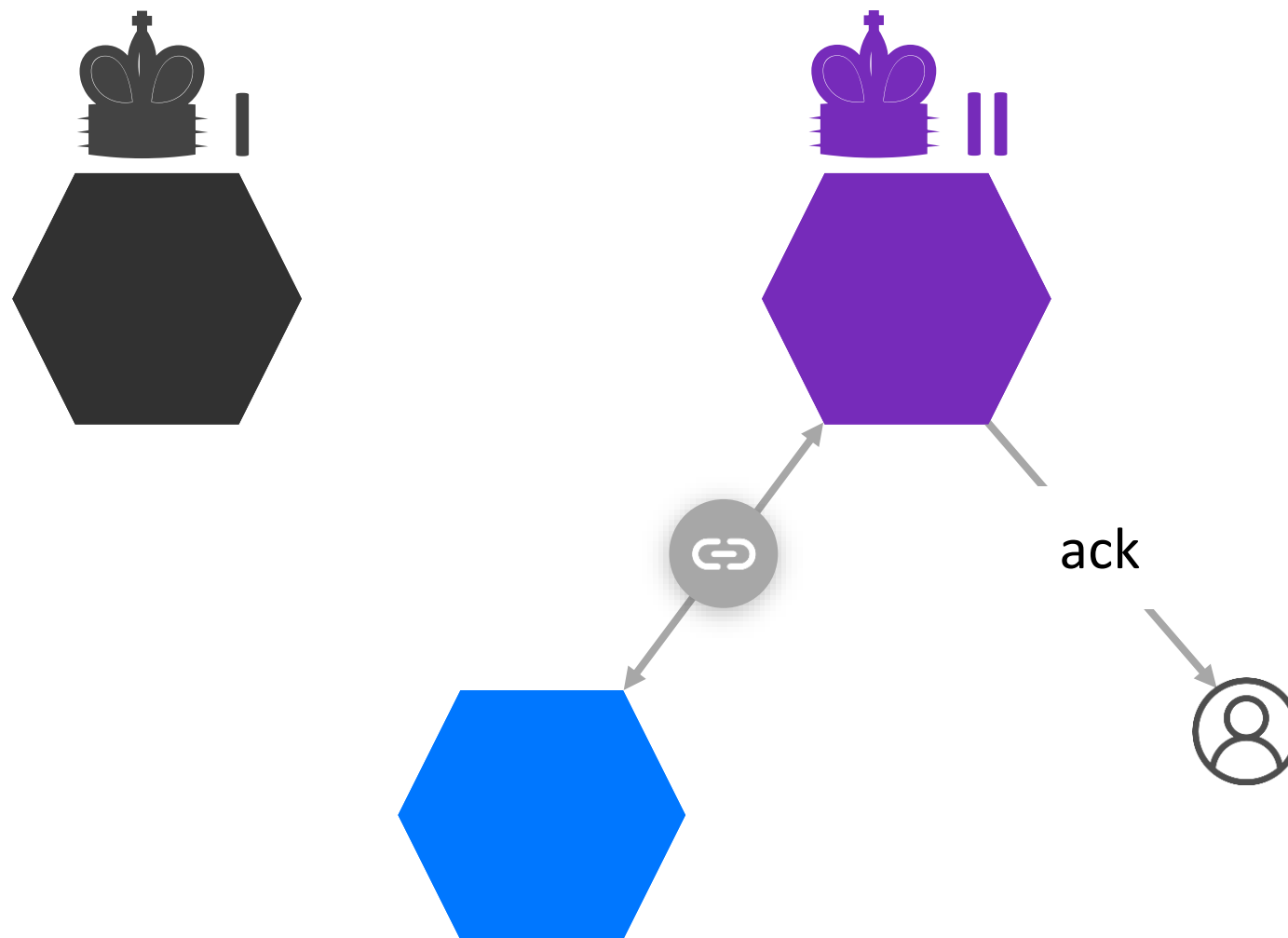
- Введение
- Особенности Raft в Tarantool
- Проблема асинхронных транзакций
- Размещение кластера Raft в двух ЦОДах
- Обнаружение различных историй лидерства
- Выводы

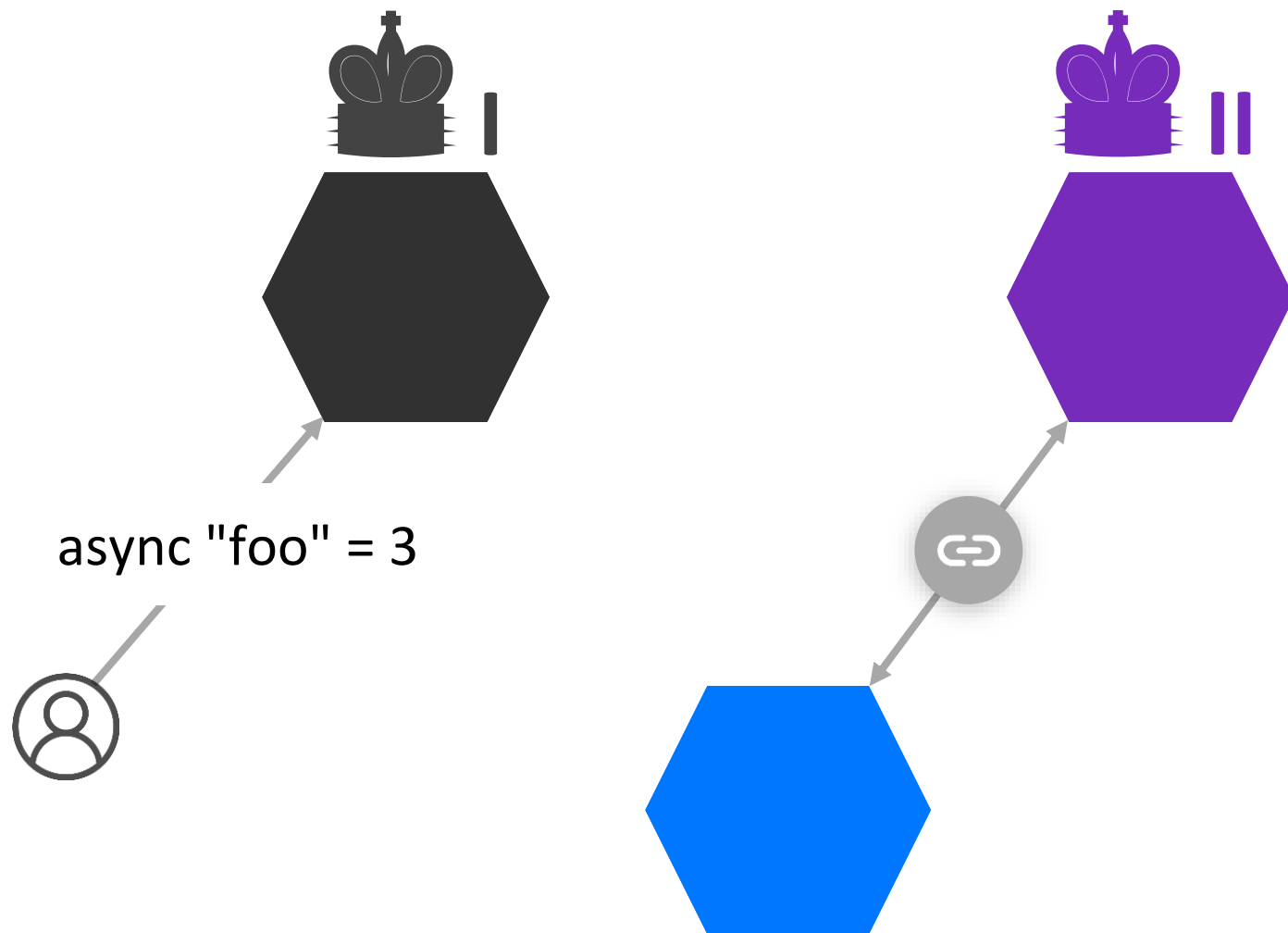


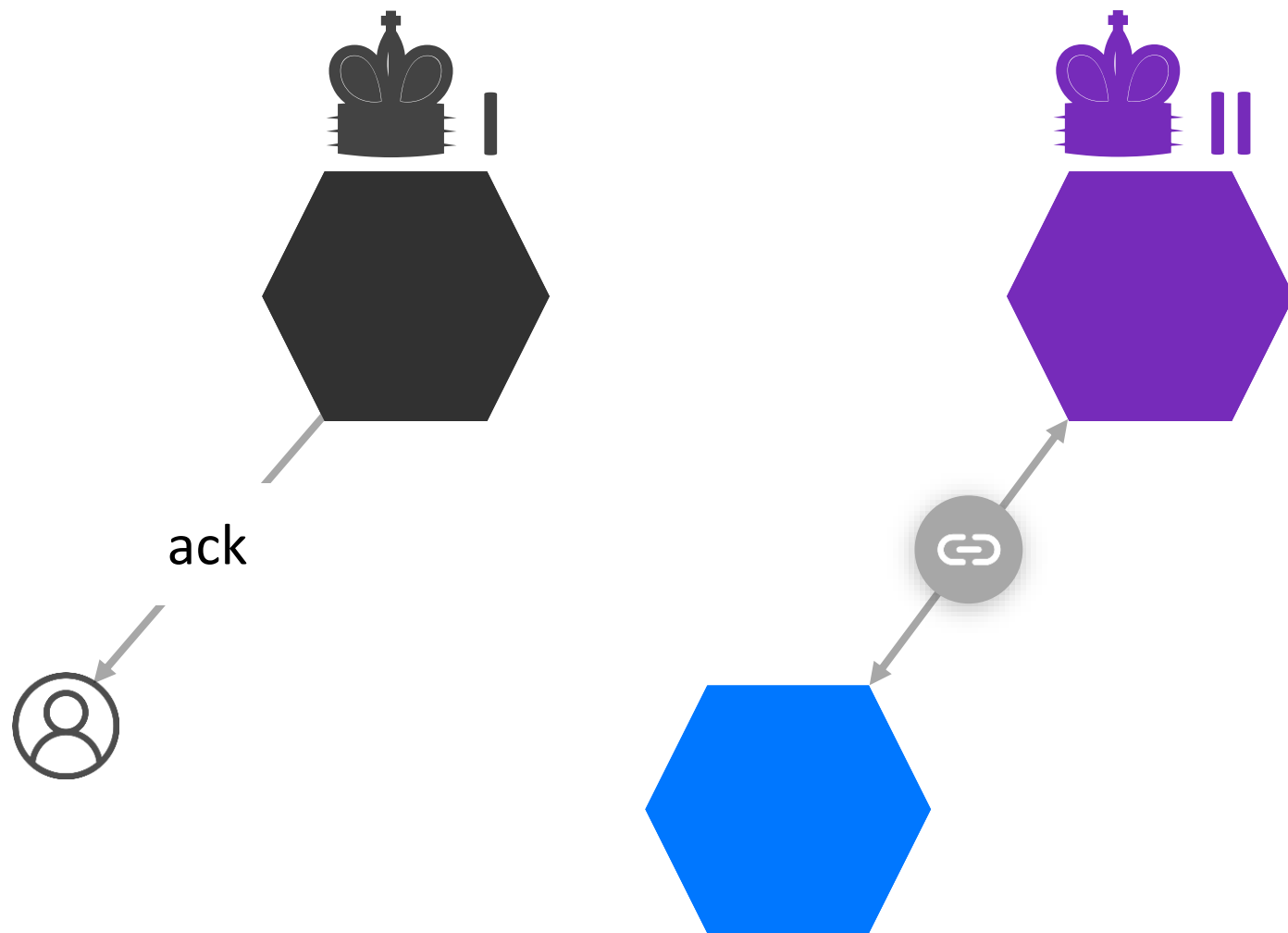


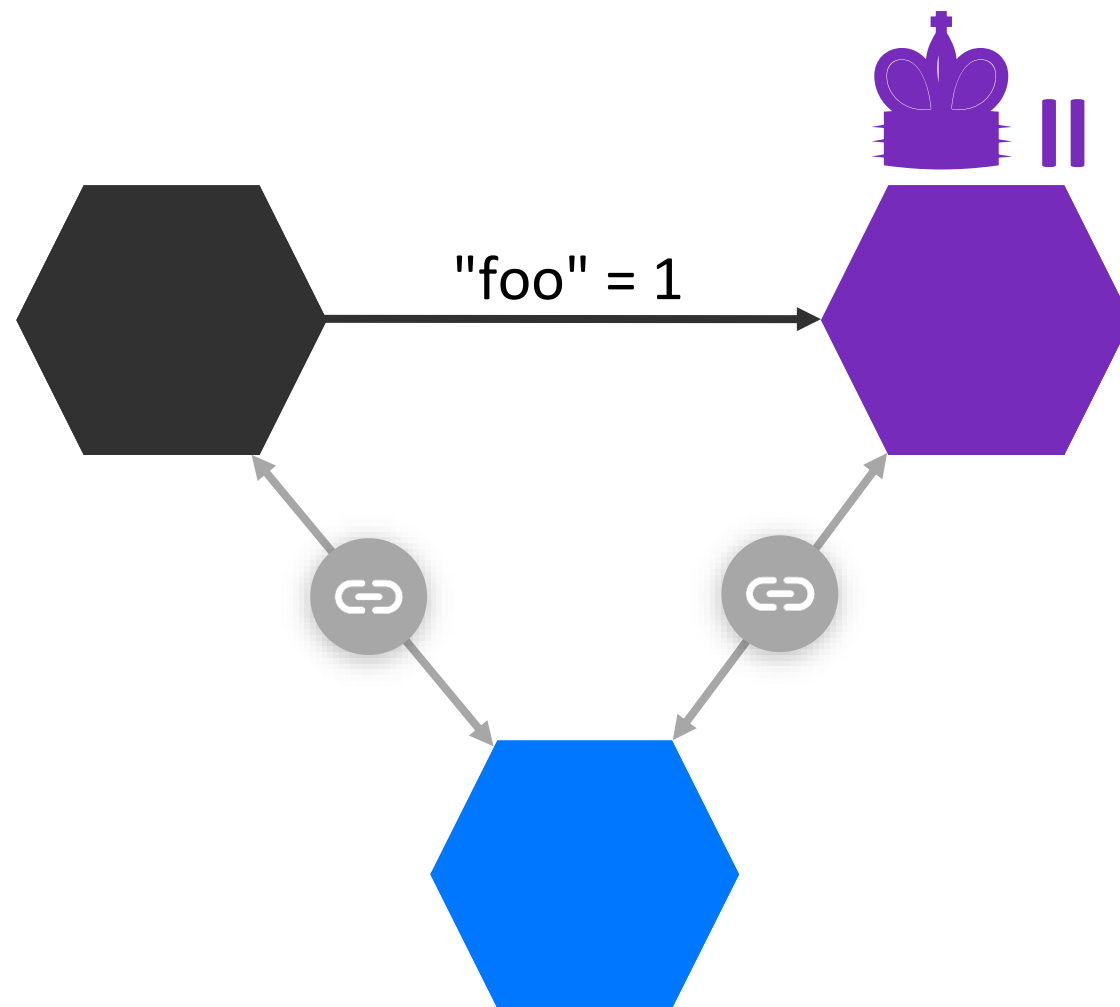


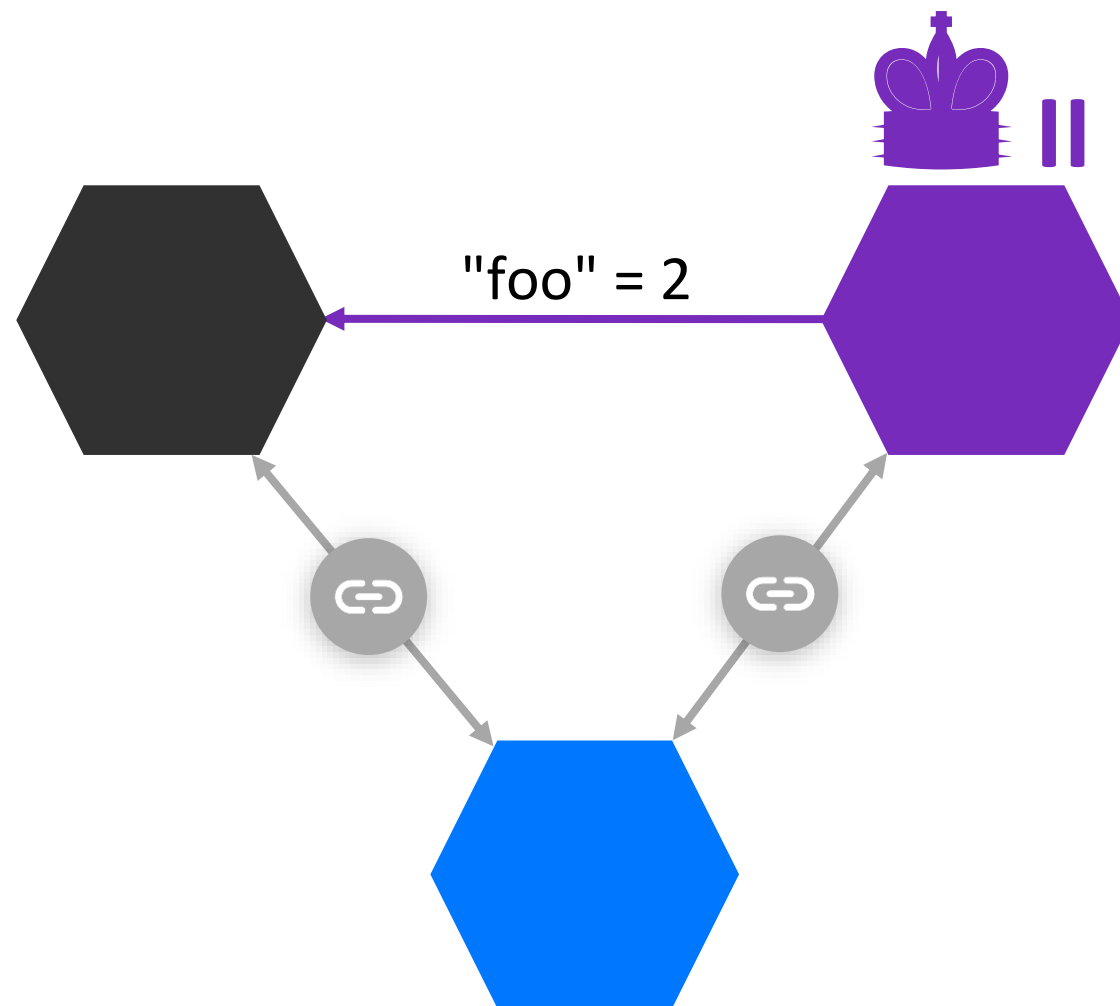


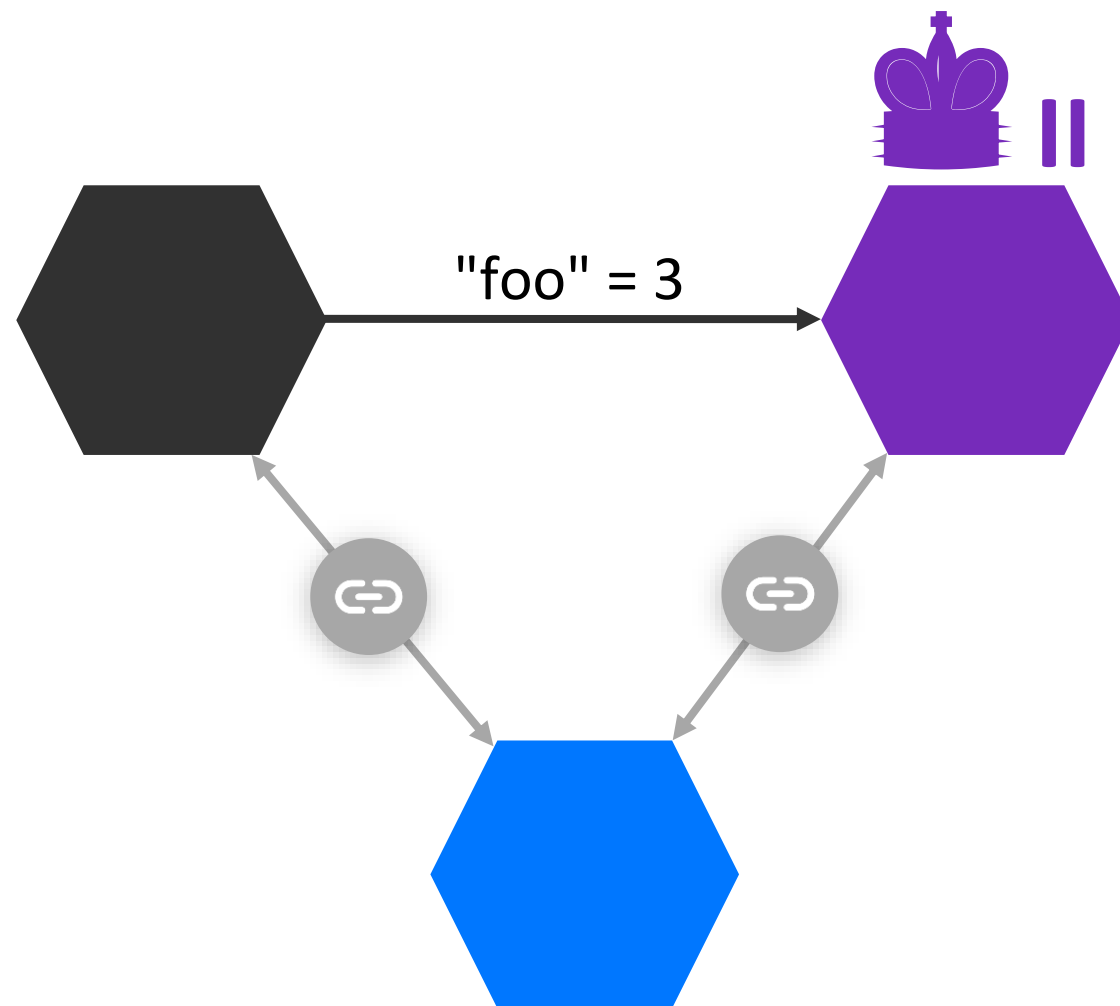


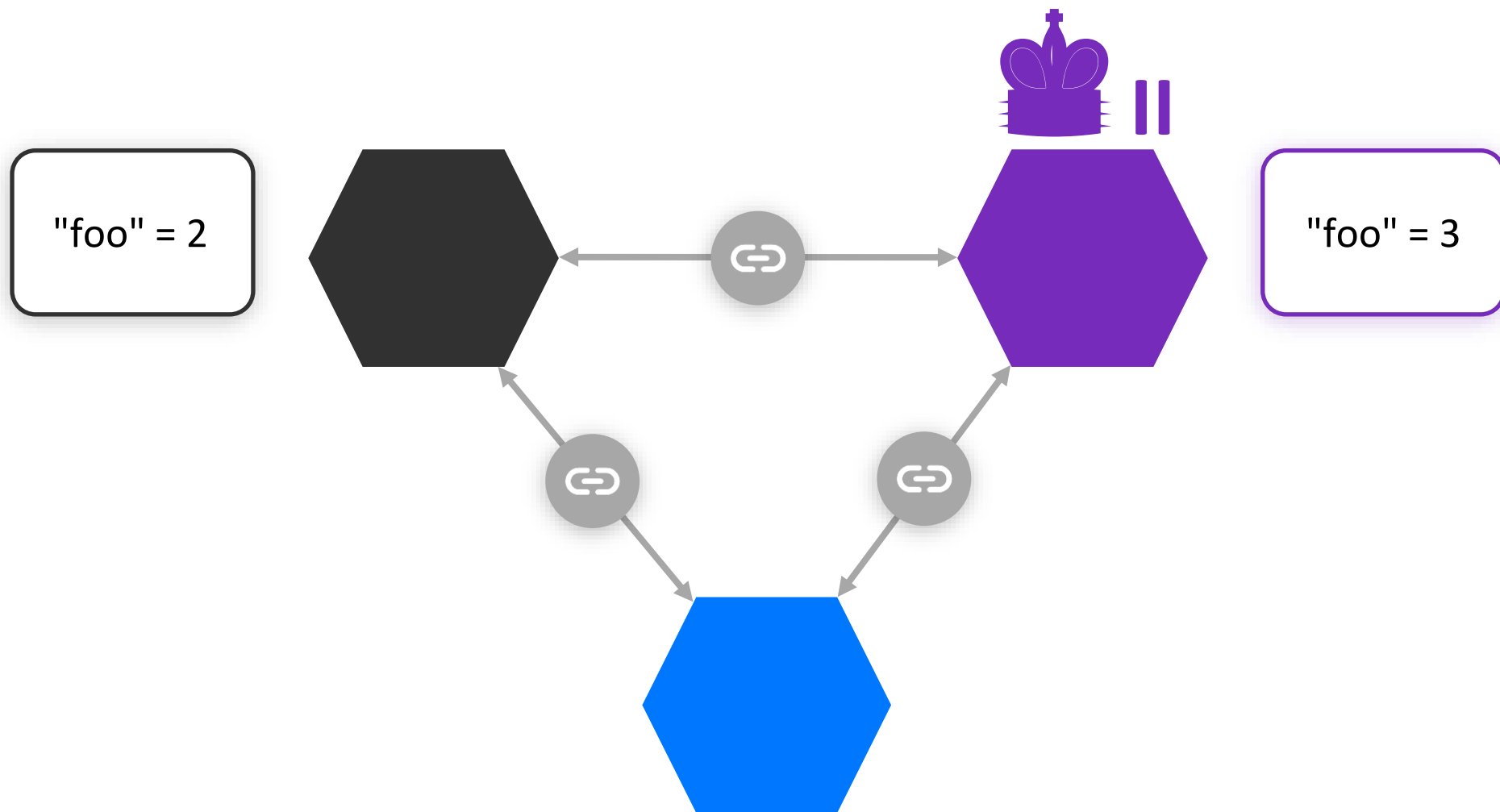












А как же CheckQuorum?

План

- Введение
- Особенности Raft в Tarantool
- Проблема асинхронных транзакций
- Размещение кластера Raft в двух ЦОДах
- Обнаружение различных историй лидерства
- Выводы

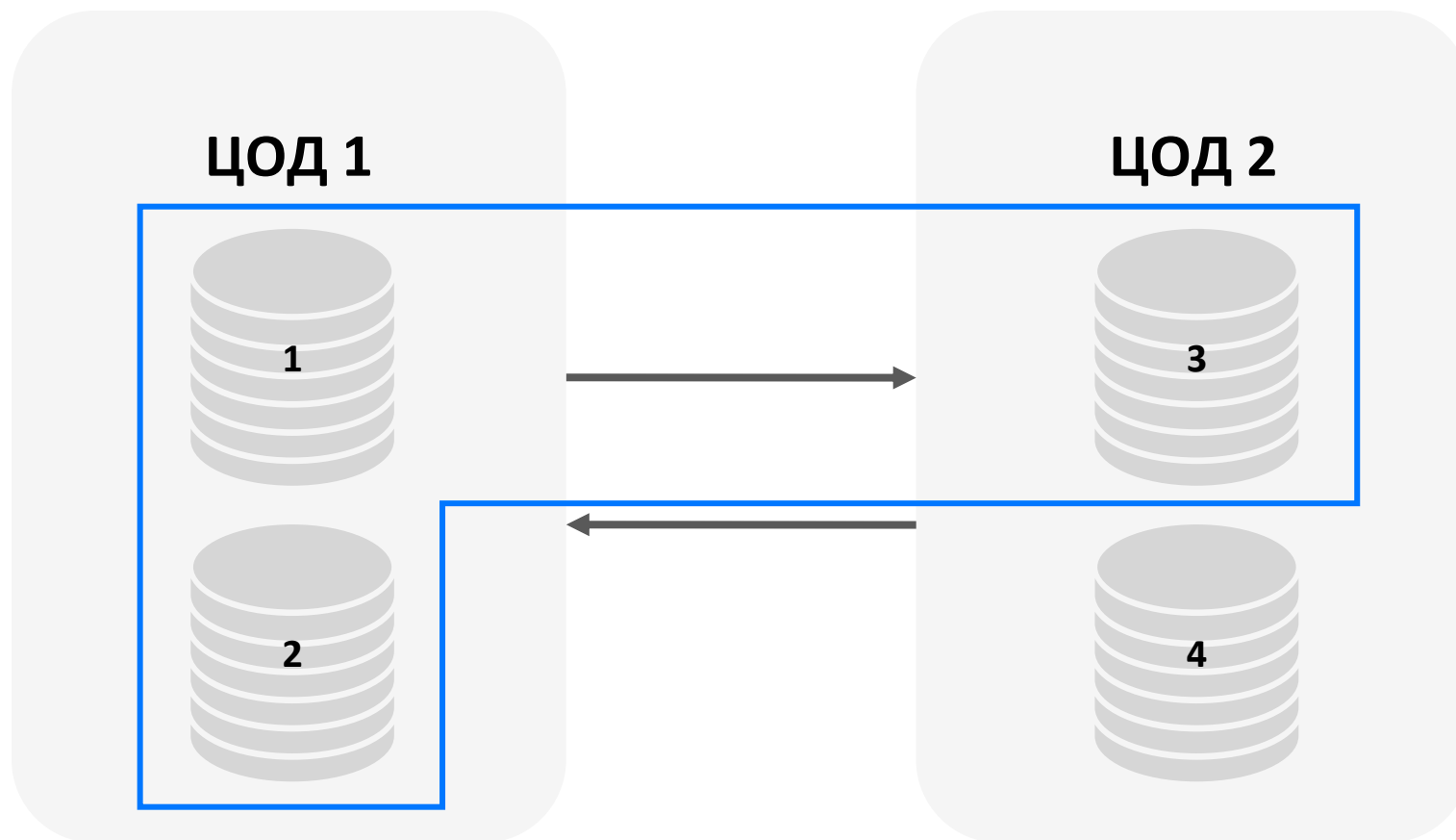
Варианты размещения кластера в двух ЦОДах

- Поровну

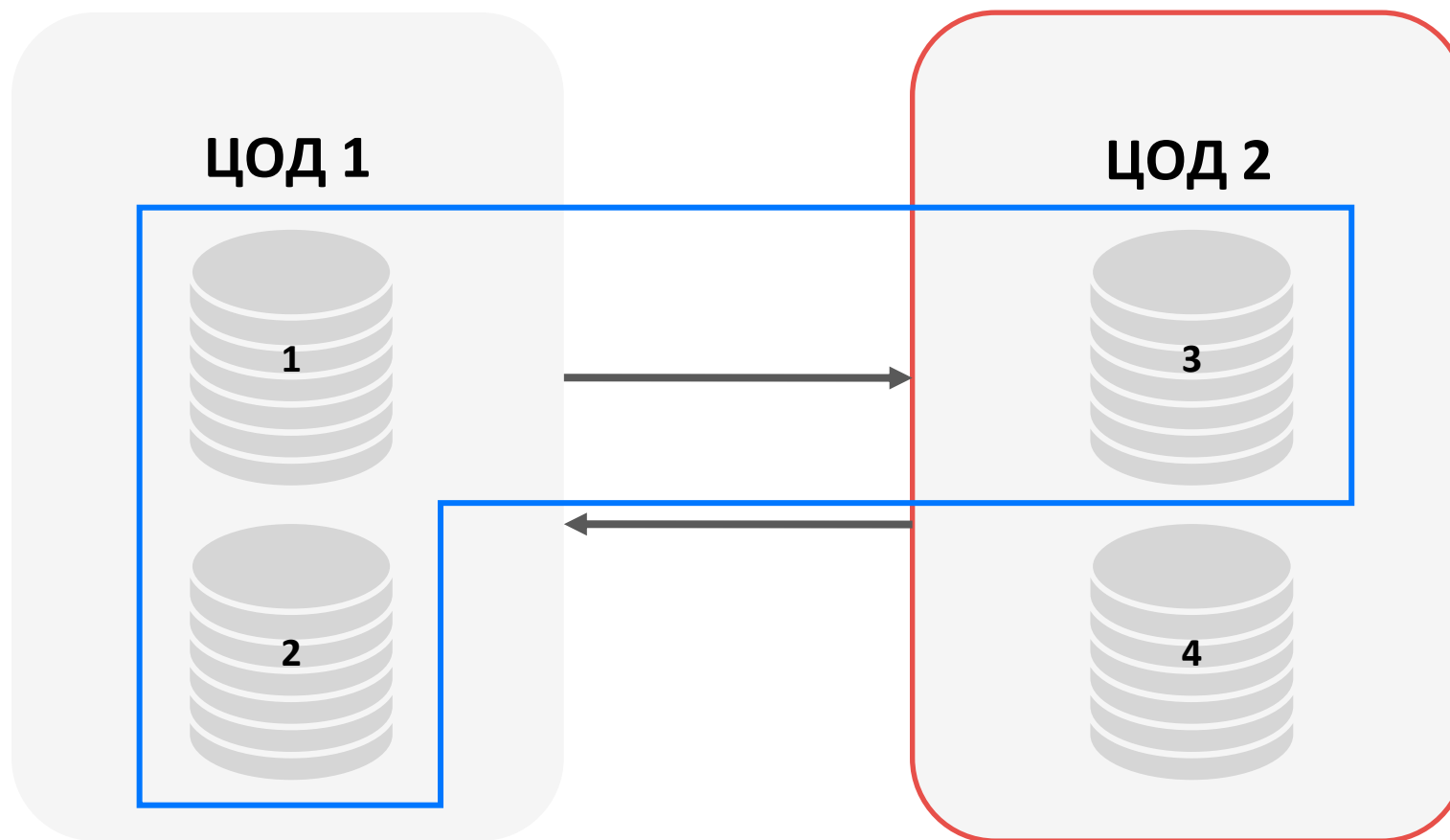
Варианты размещения кластера в двух ЦОДах

- Поровну
- Дополнительный ЦОД для голосования

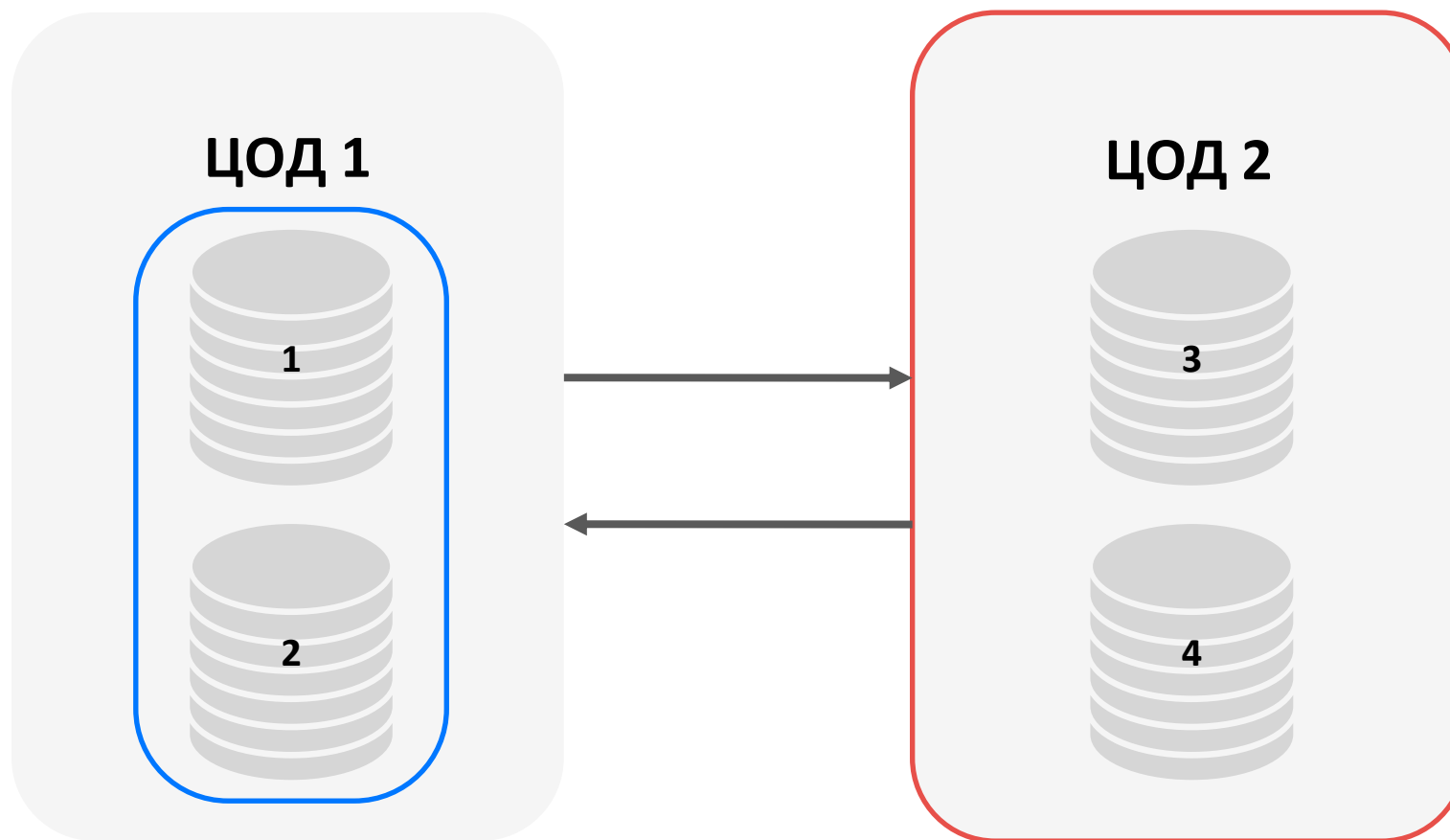
Снижение кворума при аварии



Снижение кворума при аварии

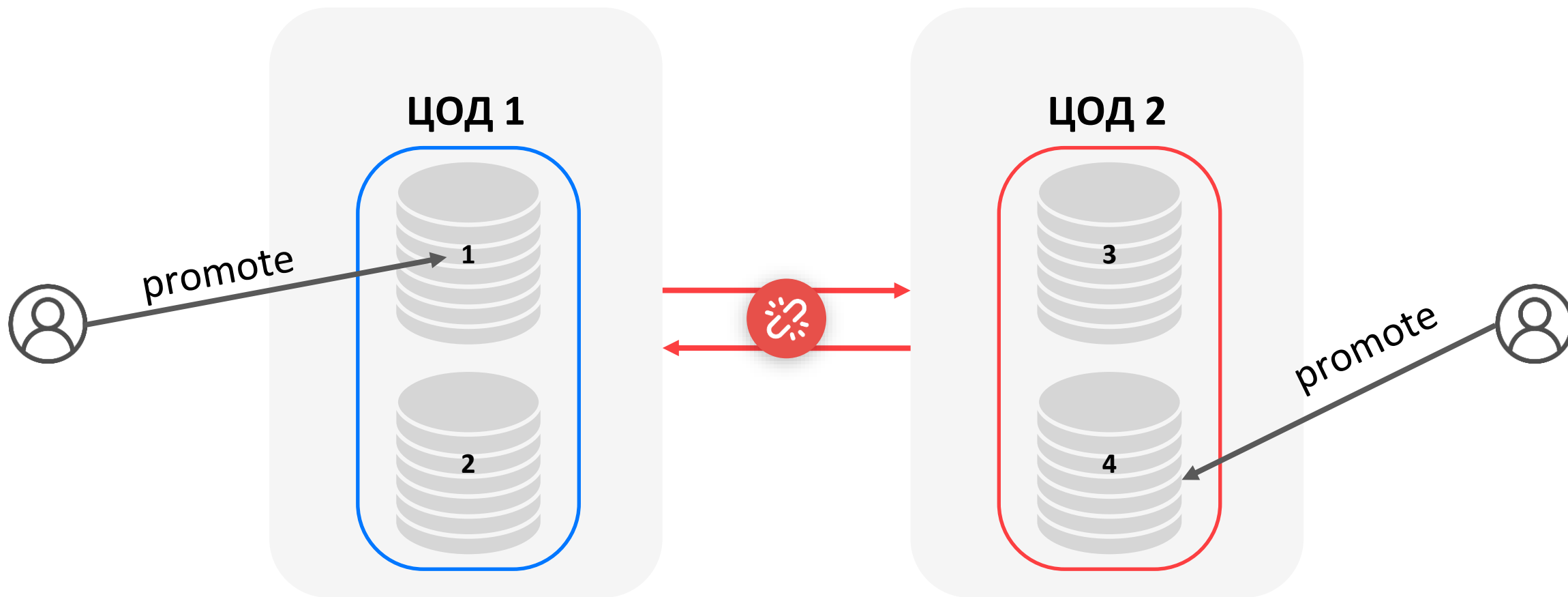


Снижение кворума при аварии

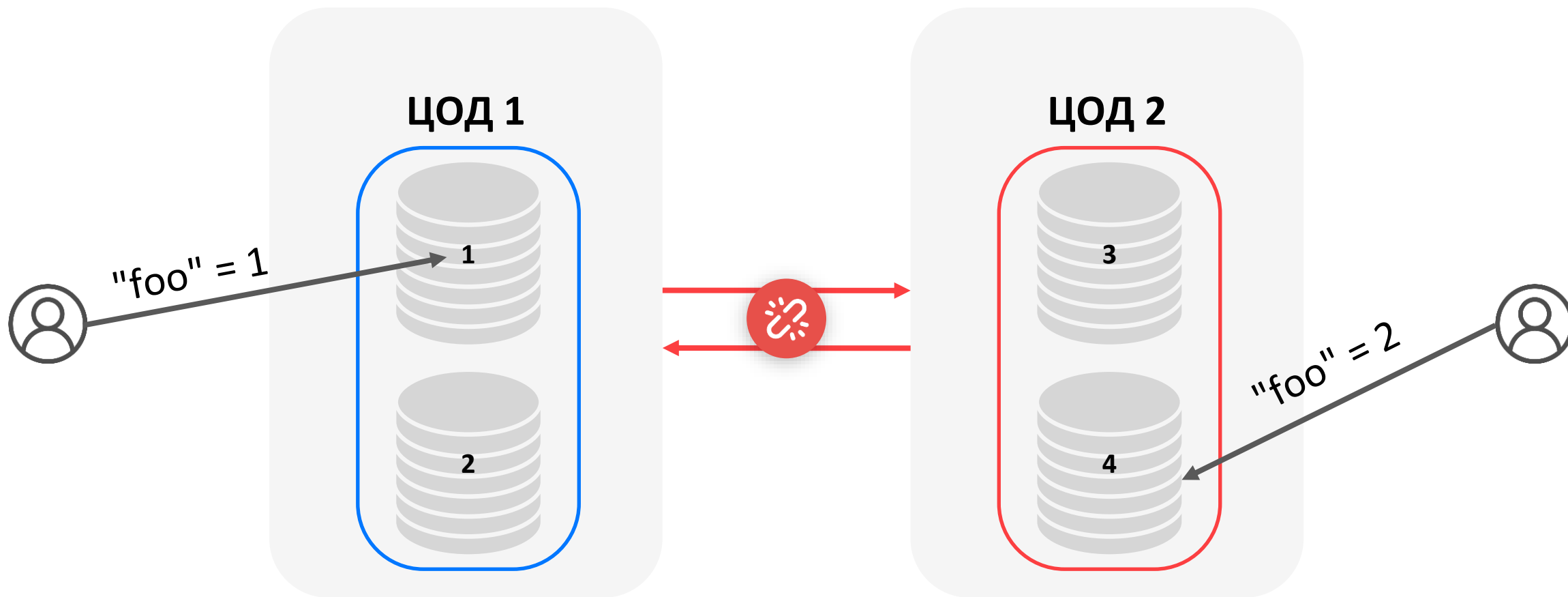


При ручном вмешательстве никто не
застрахован от ошибок,
а кластер – от сплитбрейна

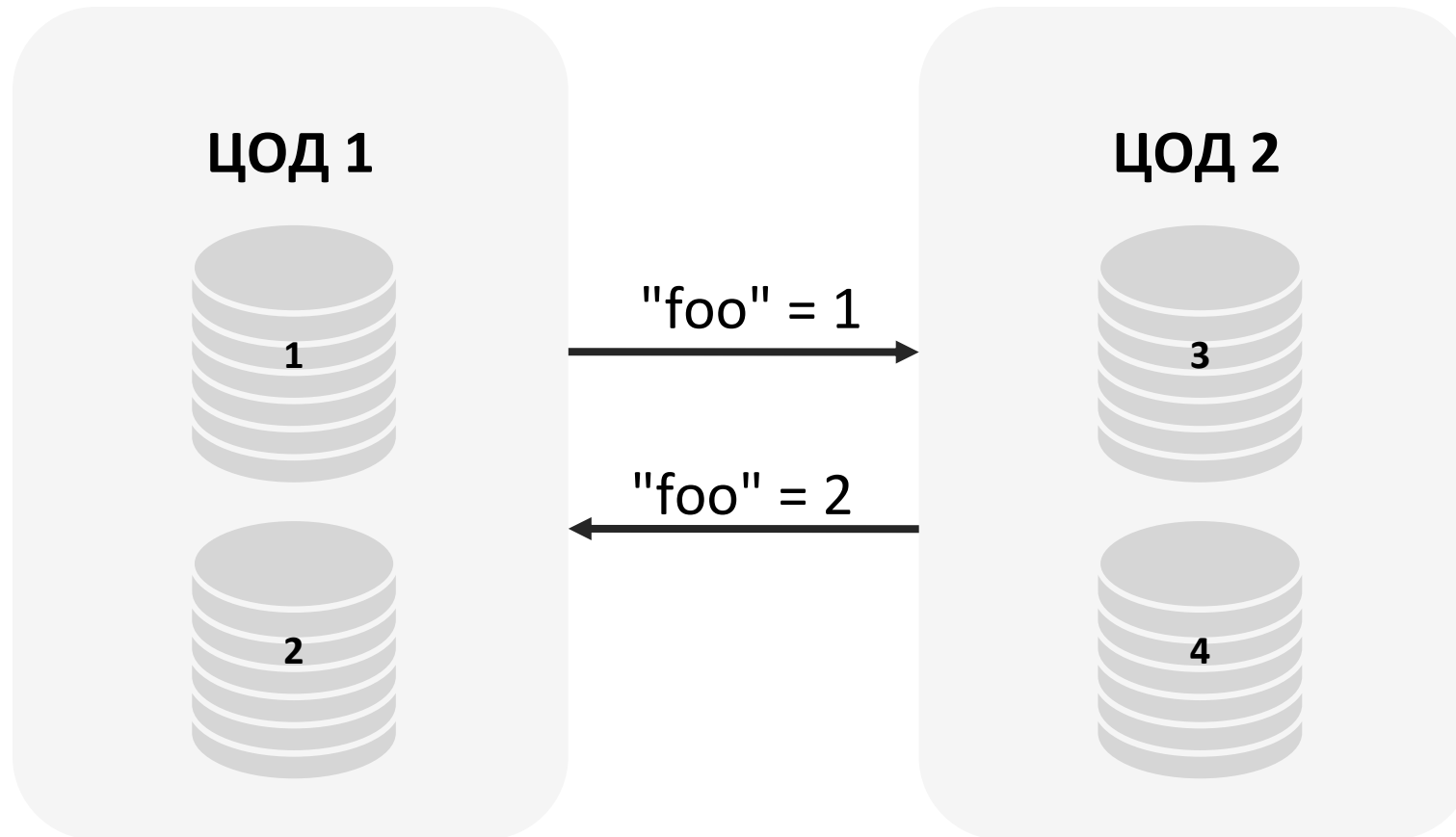
Ошибочное снижение кворума



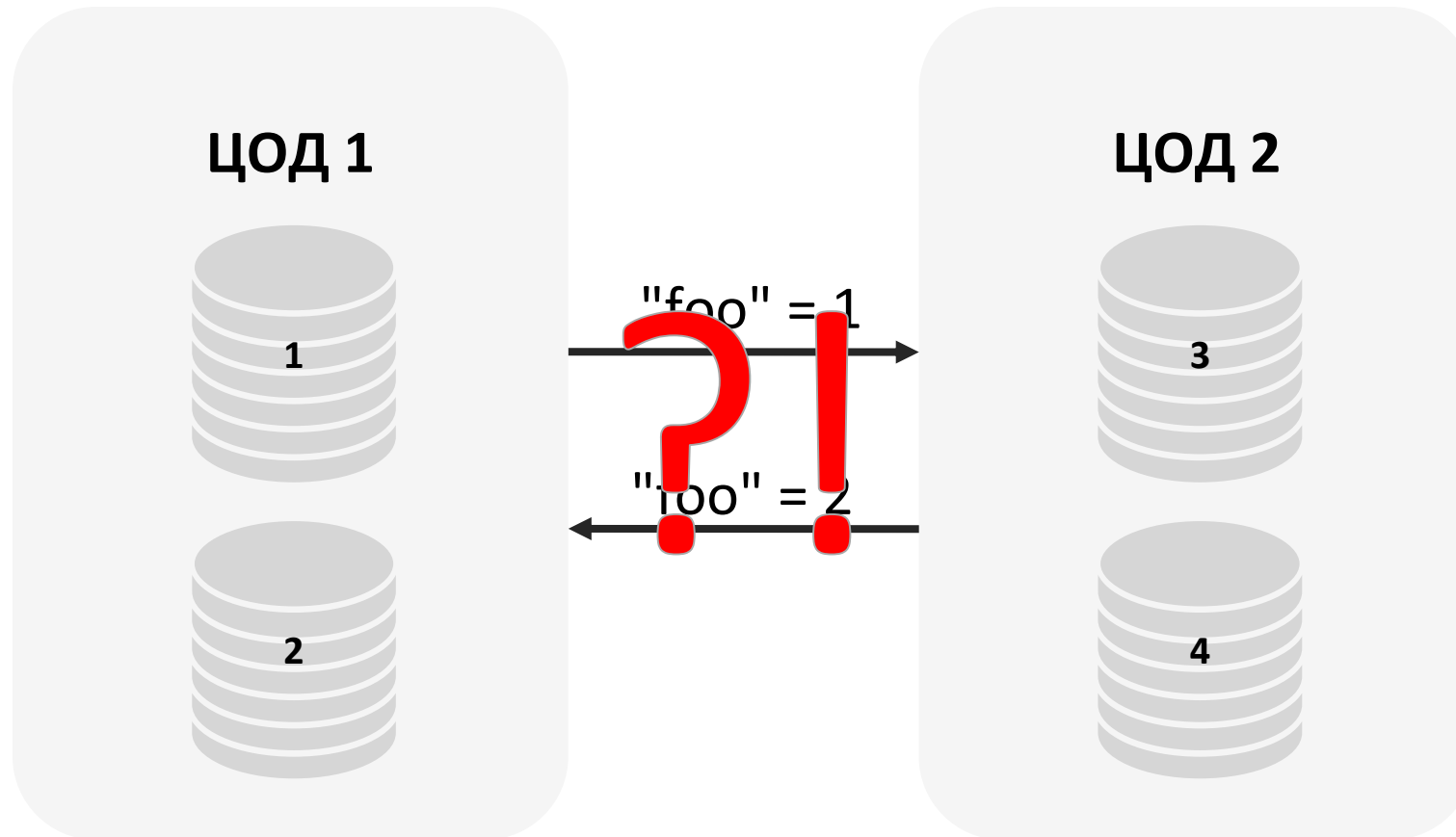
Ошибочное снижение кворума



Ошибочное снижение кворума



Ошибочное снижение кворума



Человеческие ошибки случаются и без аварий



Кворум может быть
неправильно
skonфигурирован

Человеческие ошибки случаются и без аварий



Кворум может быть
неправильно
сконфигурирован



Может появиться
два лидера

Человеческие ошибки случаются и без аварий



Кворум может быть
неправильно
skonфигурирован



Может появиться
два лидера



Нужно справляться
с последствиями работы
двух лидеров
и для синхронных
транзакций

План

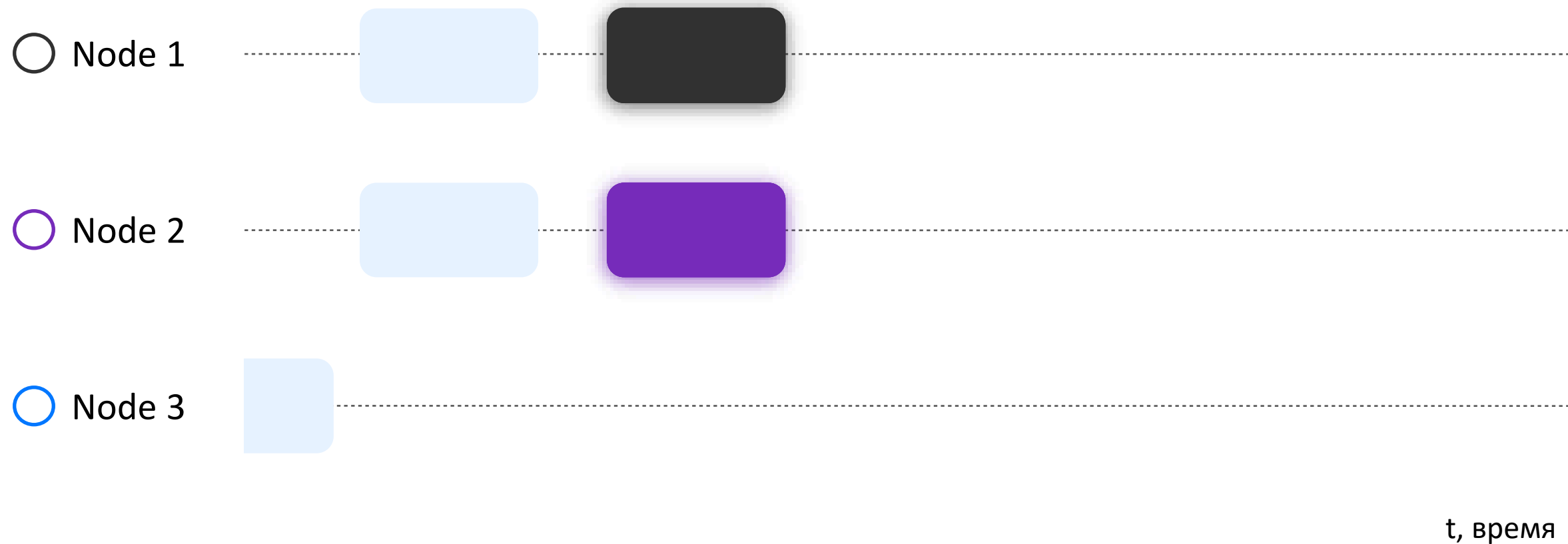
- Введение
- Особенности Raft в Tarantool
- Проблема асинхронных транзакций
- Размещение кластера Raft в двух ЦОДах
- Обнаружение различных историй лидерства
- Выводы

Репликация журнала

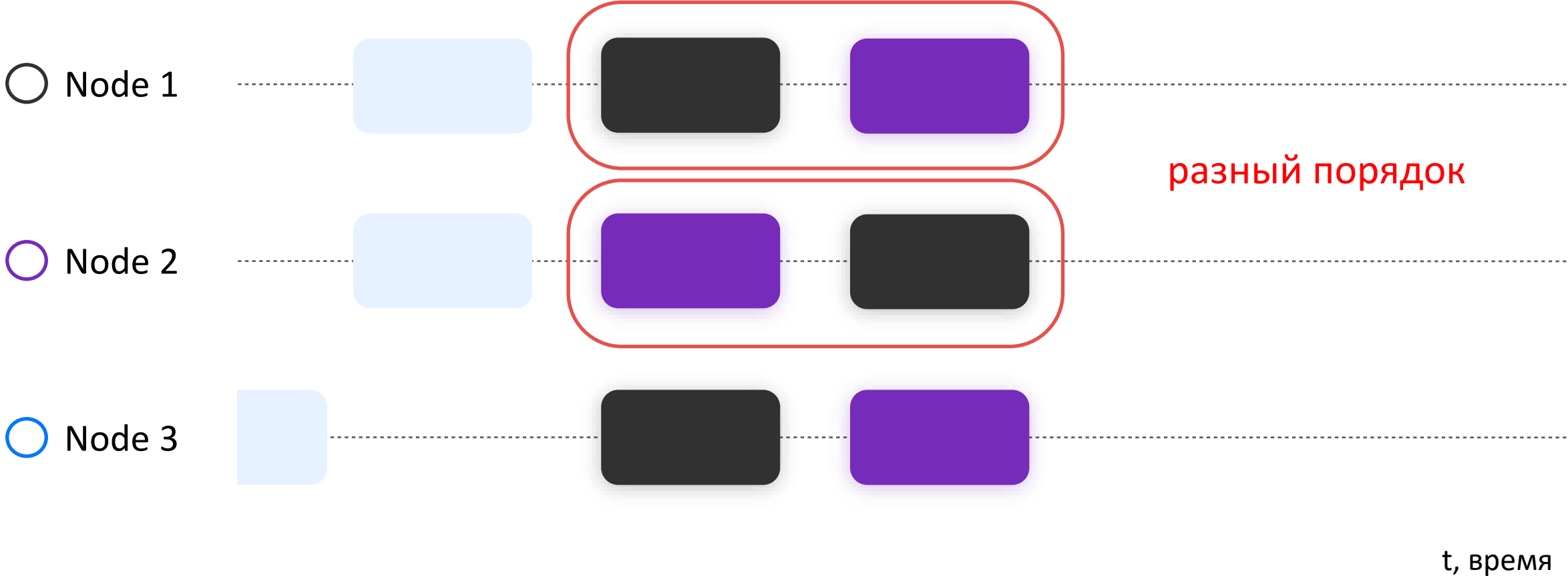


t, время

Репликация журнала



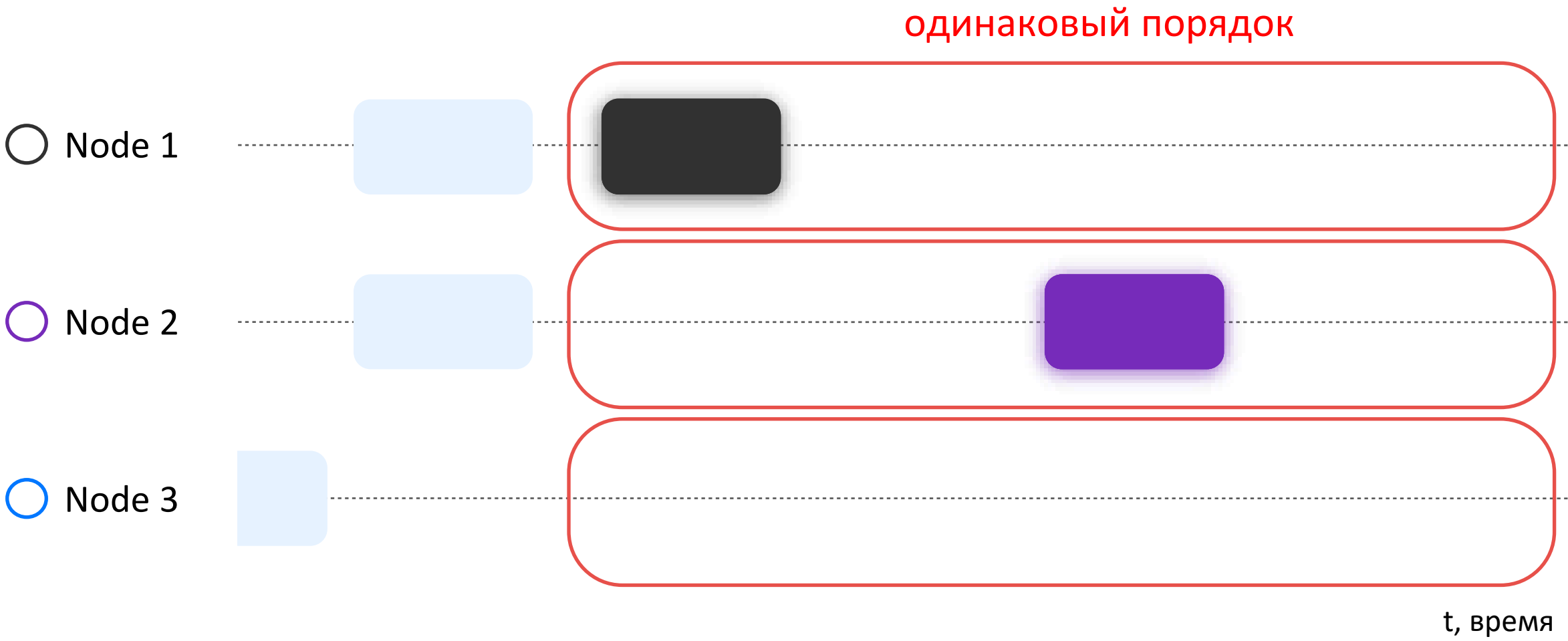
Репликация журнала



Репликация журнала с Raft



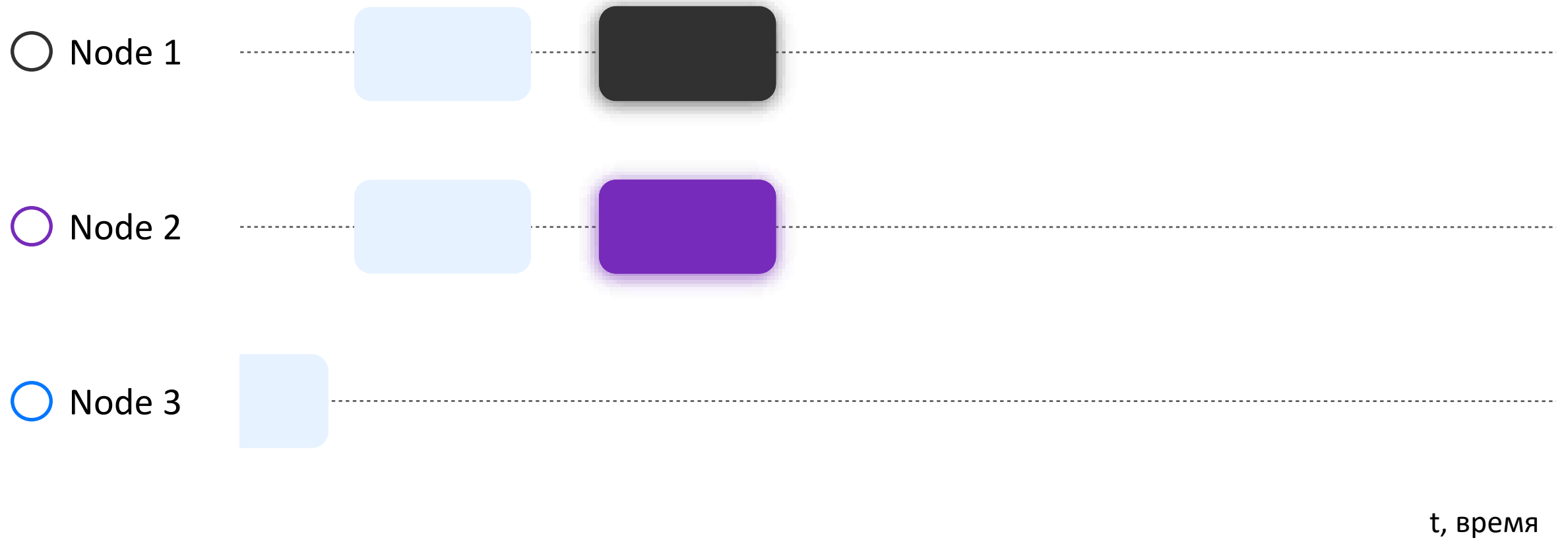
Репликация журнала с Raft



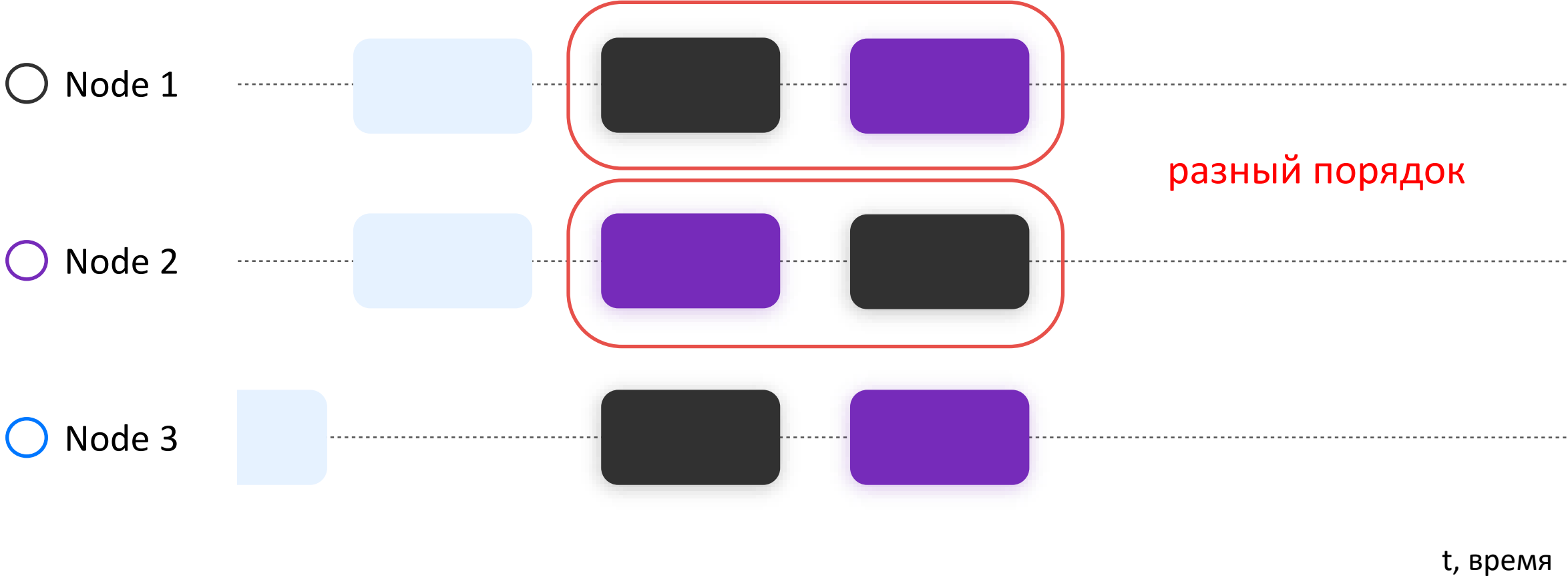
Репликация журнала с Raft. Split-brain



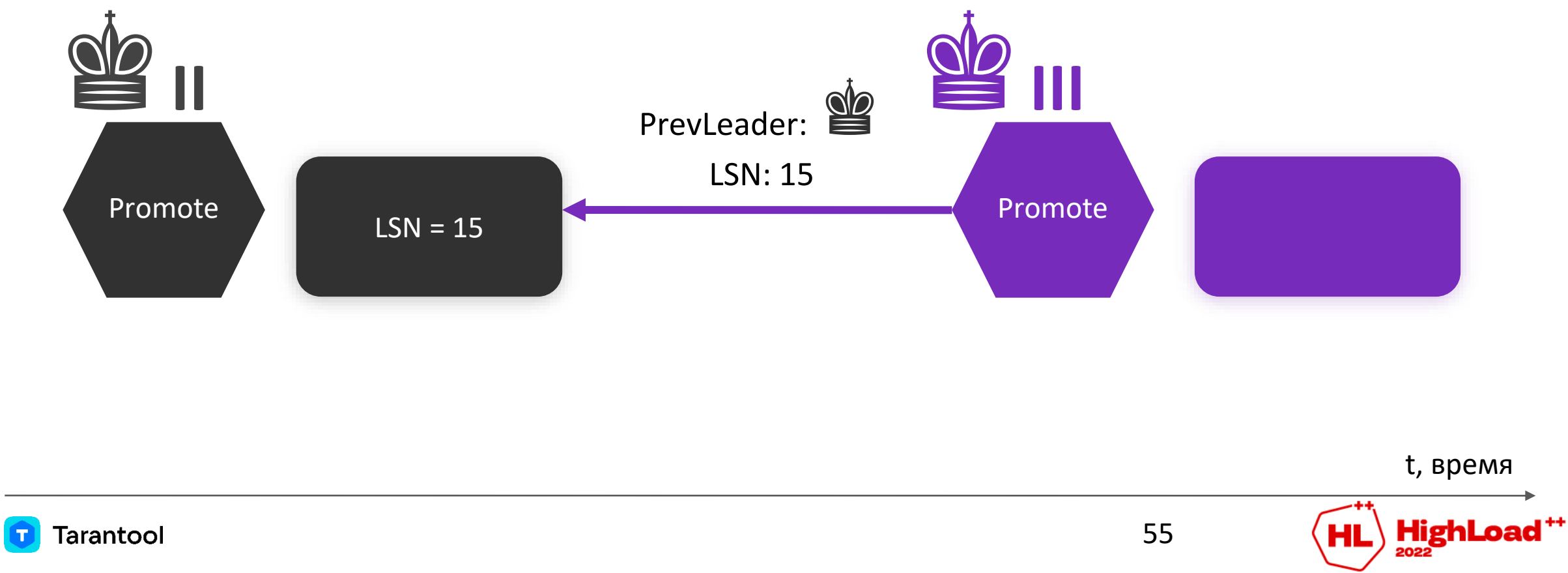
Репликация журнала с Raft. Split-brain



Репликация журнала с Raft. Split-brain

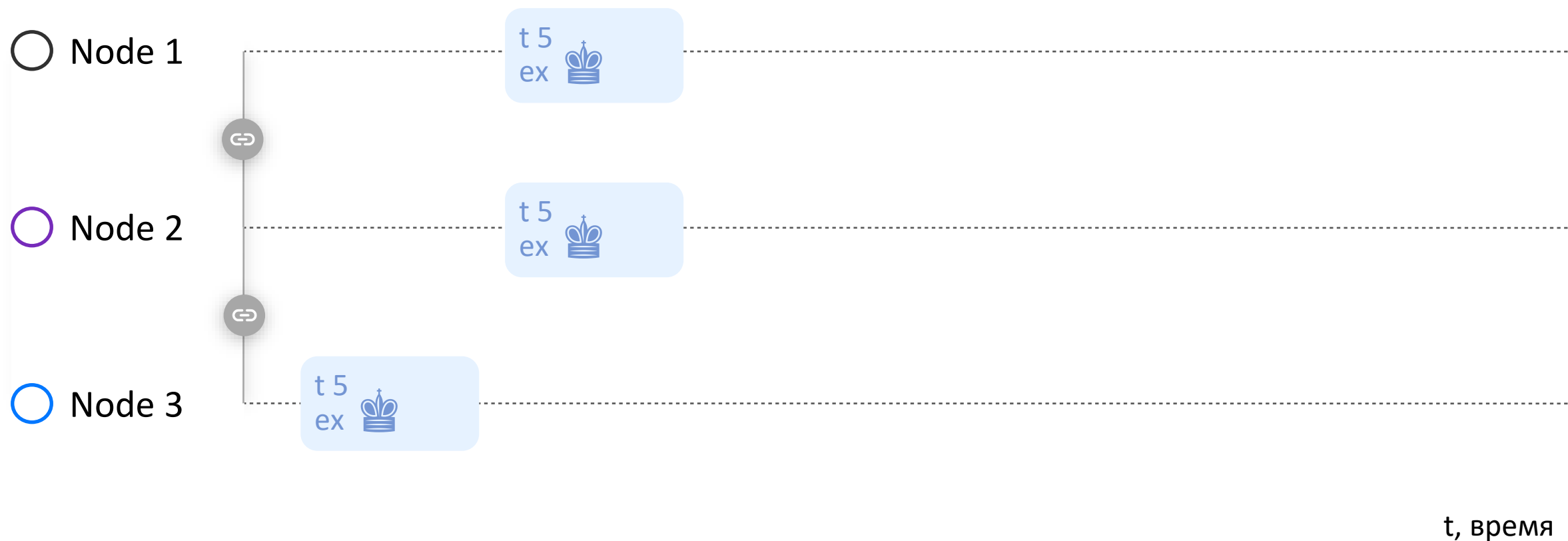


Владение журналом

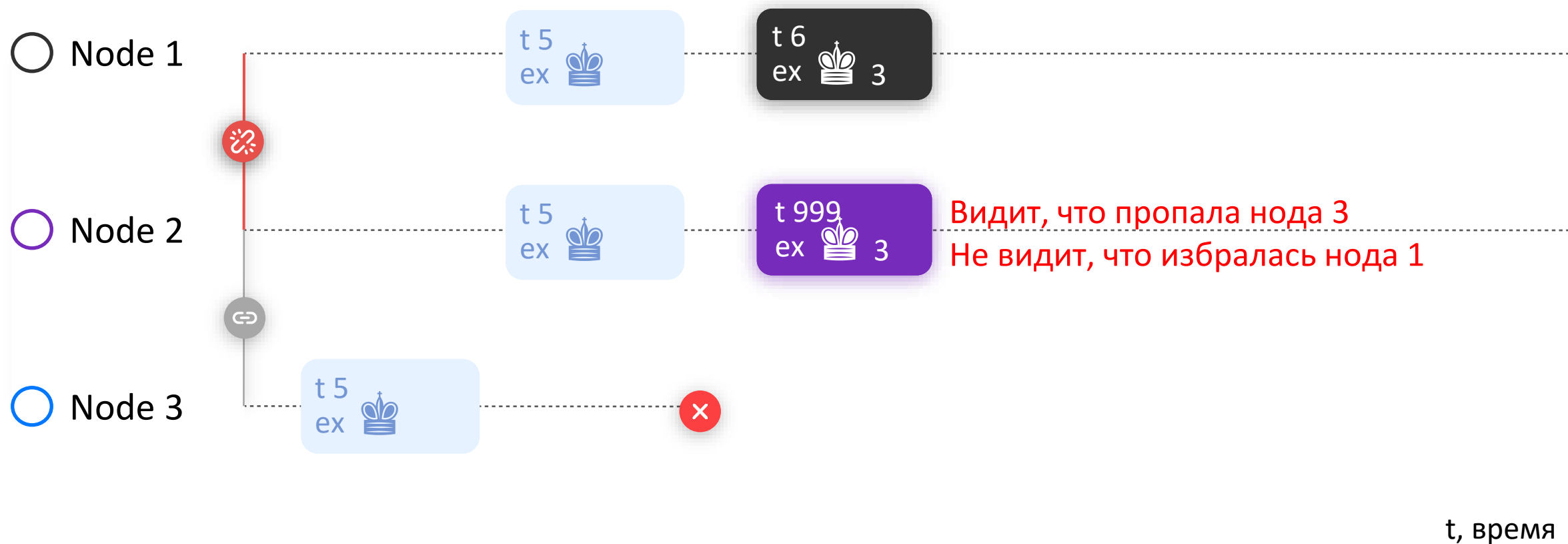


Когда всё хорошо, PROMOTE разных узлов
выстраиваются
в линейную историю

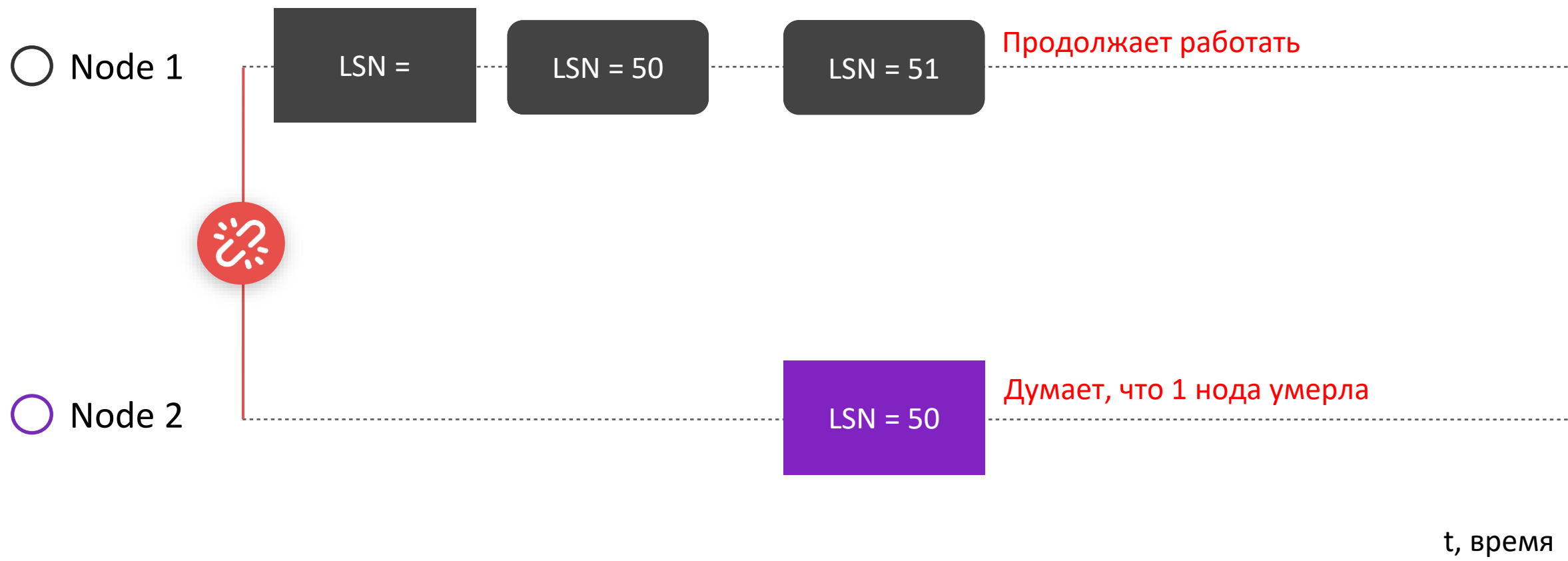
Записи в журнал при изменении лидера



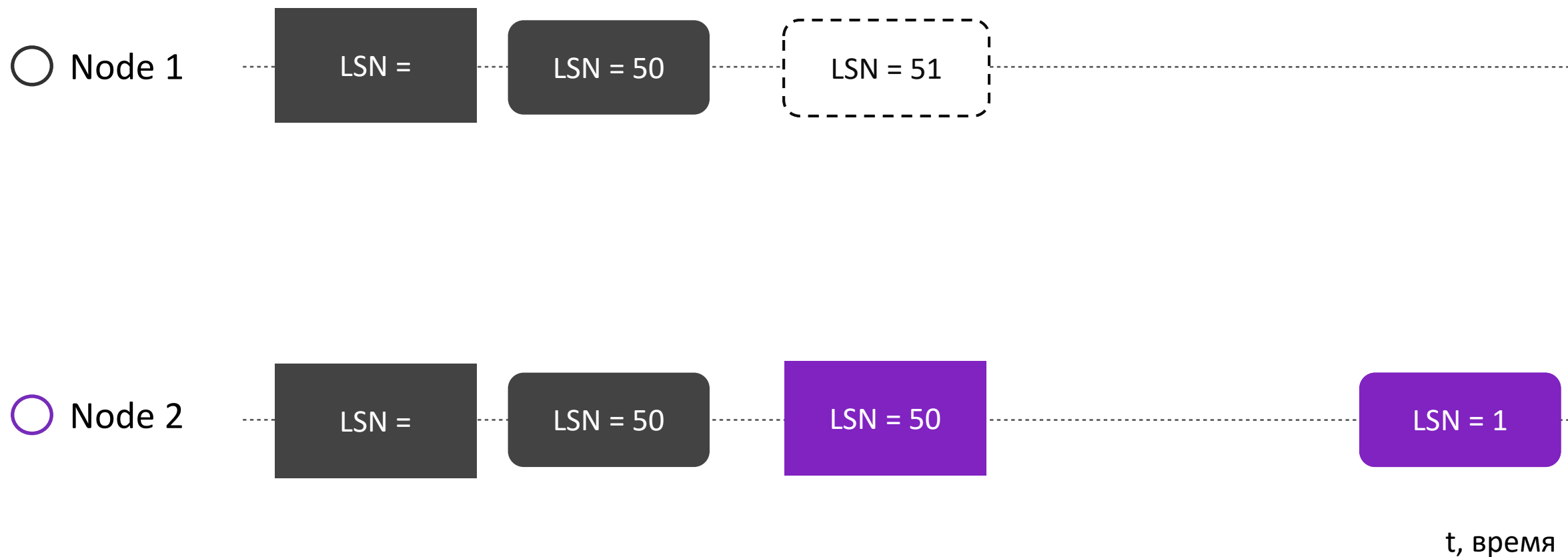
Записи в журнал при изменении лидера



Старый лидер подтвердил транзакцию



Старый лидер подтвердил транзакцию



Почему мы не разрешаем
конфликты автоматически?

План

- Введение
- Особенности Raft в Tarantool
- Проблема асинхронных транзакций
- Размещение кластера Raft в двух ЦОДах
- Обнаружение различных историй лидерства
- Выводы

Выводы

При наличии асинхронных транзакций даже Raft не защищает от перехода в master - master

Выводы

При наличии асинхронных транзакций даже Raft не защищает от перехода в master - master

Неправильная конфигурация вообще позволяет получить master-master даже с синхронной репликацией

Выводы

При наличии асинхронных транзакций даже Raft не защищает от перехода в master - master

Неправильная конфигурация вообще позволяет получить master-master даже с синхронной репликацией

Данные нужно защищать от последствий независимого изменения

Выводы

При наличии асинхронных транзакций даже Raft не защищает от перехода в master - master

Неправильная конфигурация вообще позволяет получить master-master даже с синхронной репликацией

Данные нужно защищать от последствий независимого изменения

Думаем о том, как упростить возвращение убежавшего лидера в кластер

The end

- При наличии асинхронных транзакций даже Raft не защищает от перехода в master - master
- Неправильная конфигурация вообще позволяет получить master-master даже с синхронной репликацией
- Данные нужно защищать от последствий независимого изменения
- Думаем о том, как упростить возвращение убежавшего лидера в кластер

